# Learning Task Specifications for Reinforcement Learning from Human Feedback

David Lindner

Microsoft Swiss JRC Workshop, 29.03.2022

Based on joint work with Matteo Turchetta, Sebastian Tschiatschek, Kamil Ciosek, Katja Hofmann, and Andreas Krause
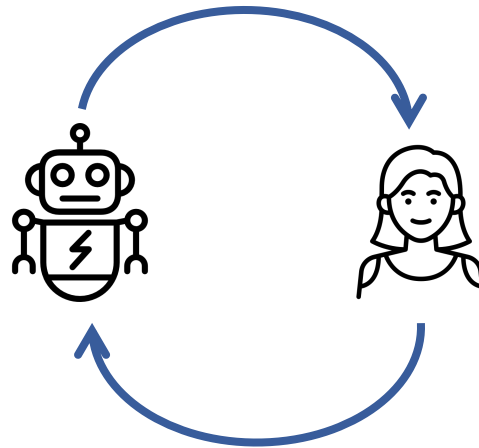
ETH zürich

LAS | Learning & Adaptive Systems

Microsoft

# **Where** do rewards come from?

Well specified
reward function

Autonomous Driving

Virtual Assistant

Reward function?
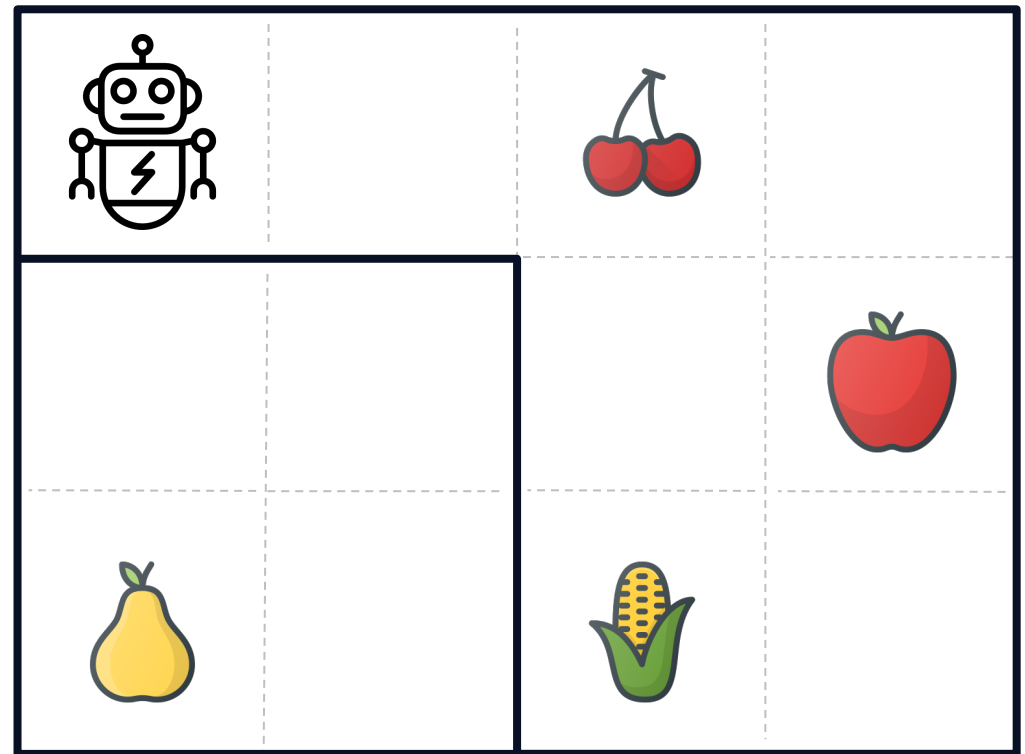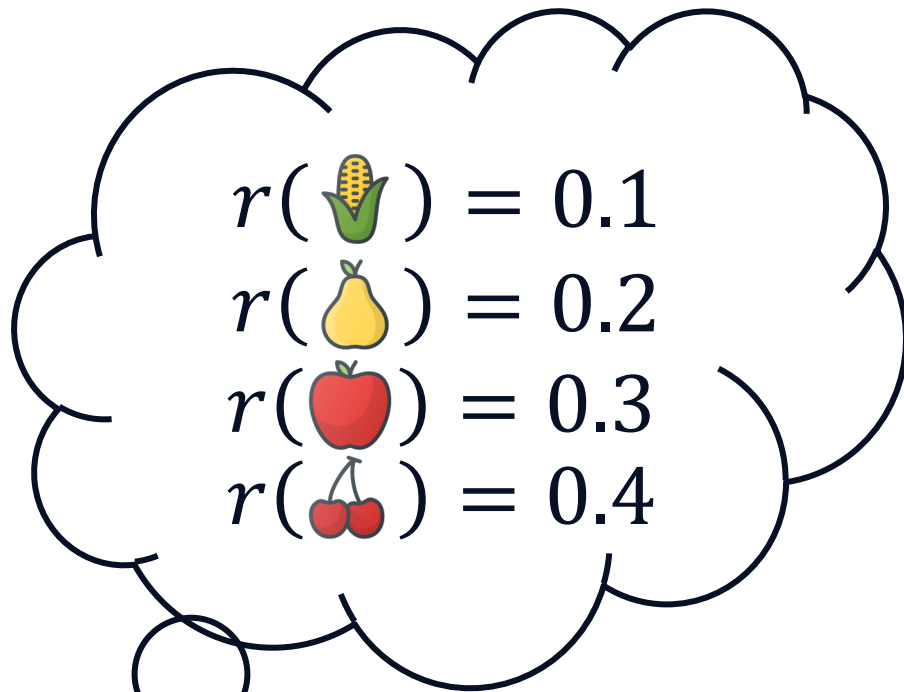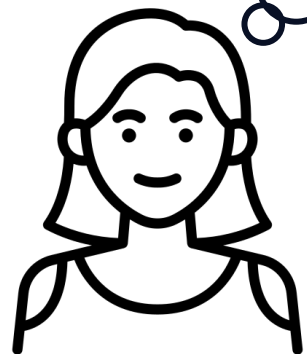
Learning from
human feedback

# How to communicate task specifications to RL agents?

Task communication

Task representation

Reward function

Goal states

Constraints

Explicit specification

Reward function

Goal states

Constraints

How should we represent tasks that humans want RL agents to do?

Implicit specification

Demonstrations

Preference labels

RL / Policy optimization

Active learning

Which queries should RL agents make to learn what humans want most effectively?

$$r(\text{🌽}) = 0.1$$
$$r(\text{🍐}) = 0.2$$
$$r(\text{🍎}) = 0.3$$
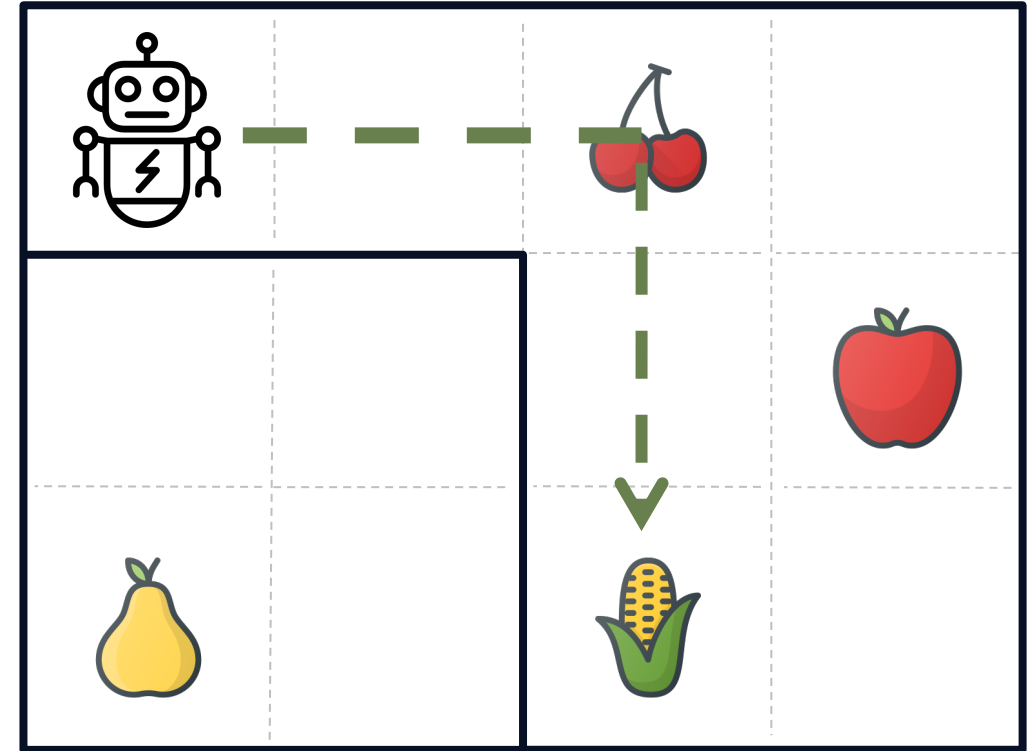$$r(\text{🍒}) = 0.4$$

$T = 4$

# Considering potentially optimal policies improves sample efficiency



- 🍒 will always be collected
- 🍐 will never be collected

🤖 only has to decide between 🍎 and 🌽.
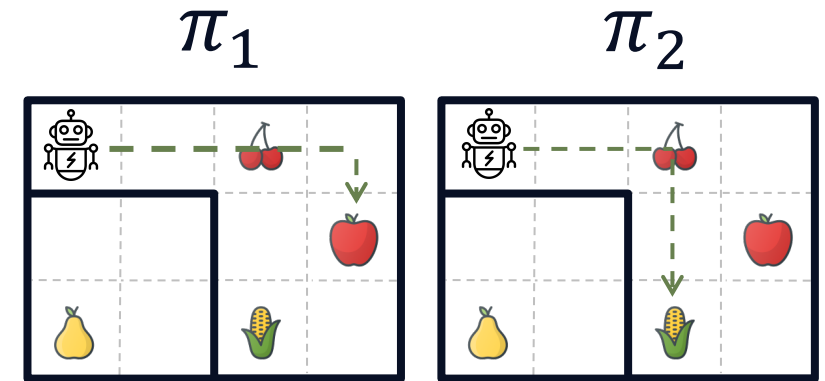
# How should we select which queries to make?

1. Consider potentially optimal policies $\pi_1$ and $\pi_2$

2. How do they differ?
Difference in return: $\hat{G}(\pi_1) - \hat{G}(\pi_2)$

3. Which observations help most to distinguish them?

$q^* \in \text{argmax}_{q \in \mathcal{Q}_c} I\big(\hat{G}(\pi_1) - \hat{G}(\pi_2); (q, \hat{y}) \big| \mathcal{D}\big)$



$\pi_1$ $\qquad$ $\pi_2$

$\hat{G}(\pi_1) = \hat{r}(🍒) + \hat{r}(🍎)$

$\hat{G}(\pi_2) = \hat{r}(🍒) + \hat{r}(🌽)$

$\hat{G}(\pi_1) - \hat{G}(\pi_1) = \hat{r}(🍎) - \hat{r}(🌽)$

Relevant states:
🍎 and 🌽

# Information directed reward learning (IDRL)

Initialize reward model

if not done

Select potentially optimal policies $\pi_1$ and $\pi_2$

Select query:
$$q^* \in \mathrm{argmax}_{q \in \mathcal{Q}_c} I\big(\widehat{G}(\pi_1) - \widehat{G}(\pi_2); (q, \hat{y})\big|\mathcal{D}\big)$$

Make query $q^*$
Update reward model

if done

Find and return policy for mean estimate of reward model

Gaussian Process (GP) reward model

$P(\hat{r}(🍎)|\mathcal{D})$

$\hat{r}(🍎)$
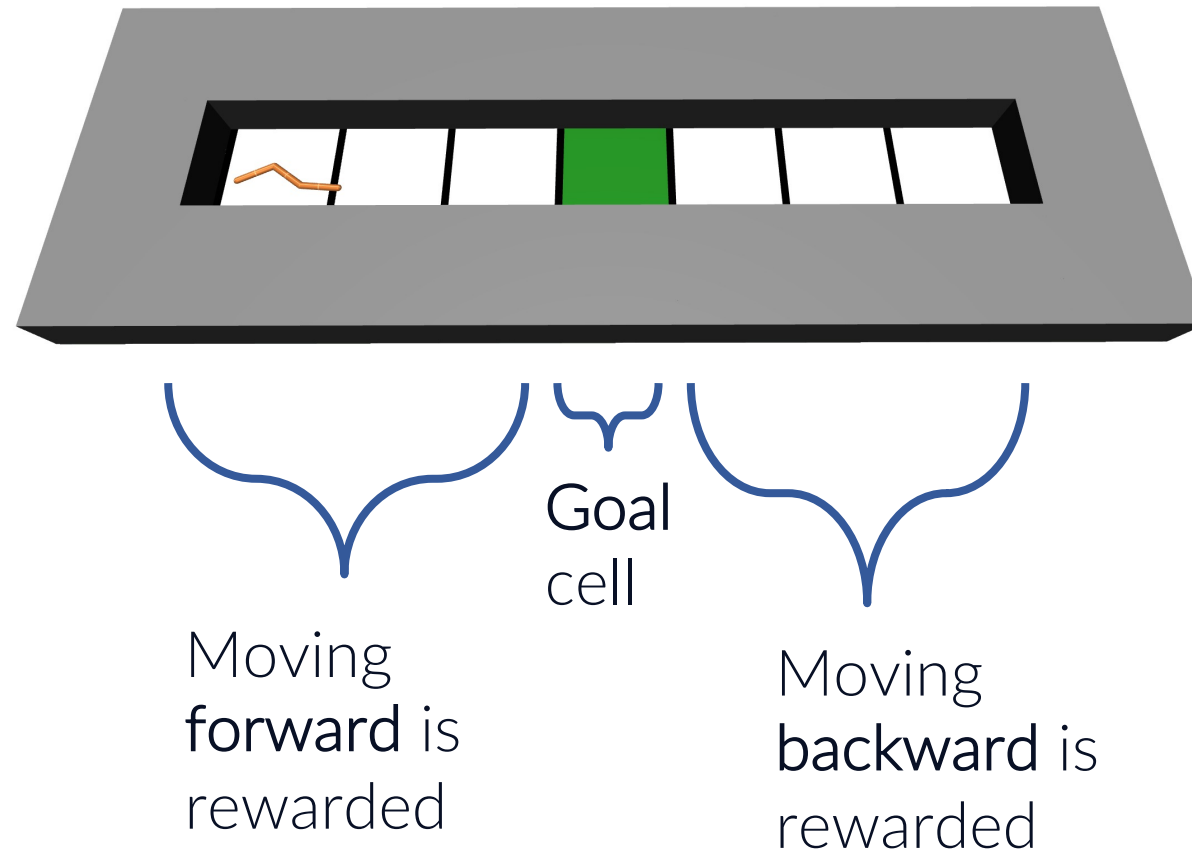
Deep Neural Network reward

- For a GP model we can compute the information gain exactly and efficiently

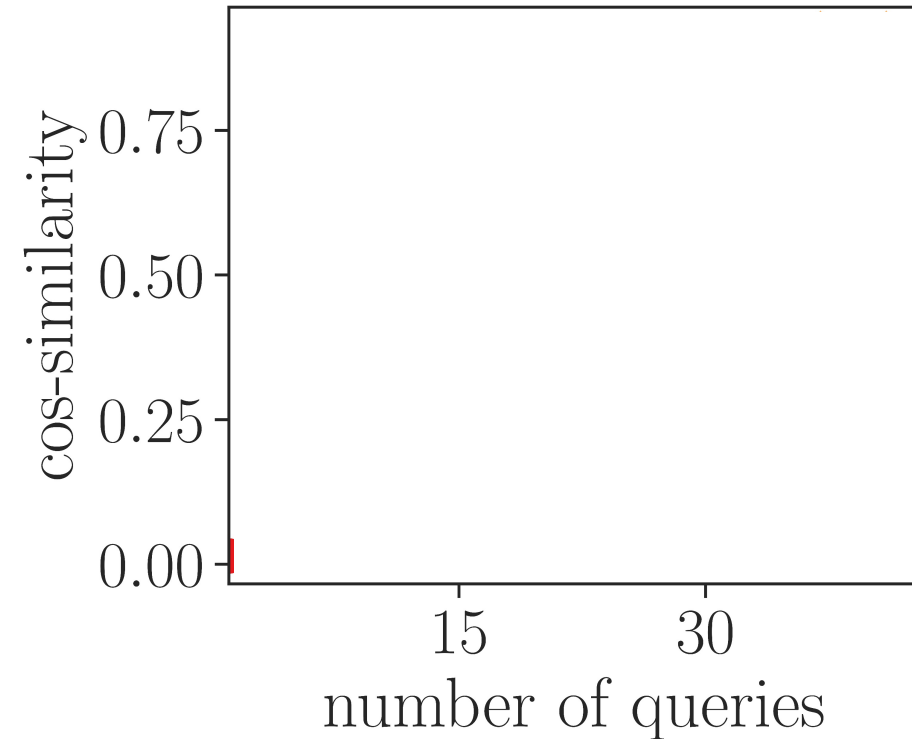- For some DNN models, we **can approximate the IDRL objective efficiently**

# IDRL can learn a complex task in the MuJoCo simulator from numerical evaluations of clips of trajectories

Goal cell

Moving **forward** is rewarded

Moving **backward** is rewarded

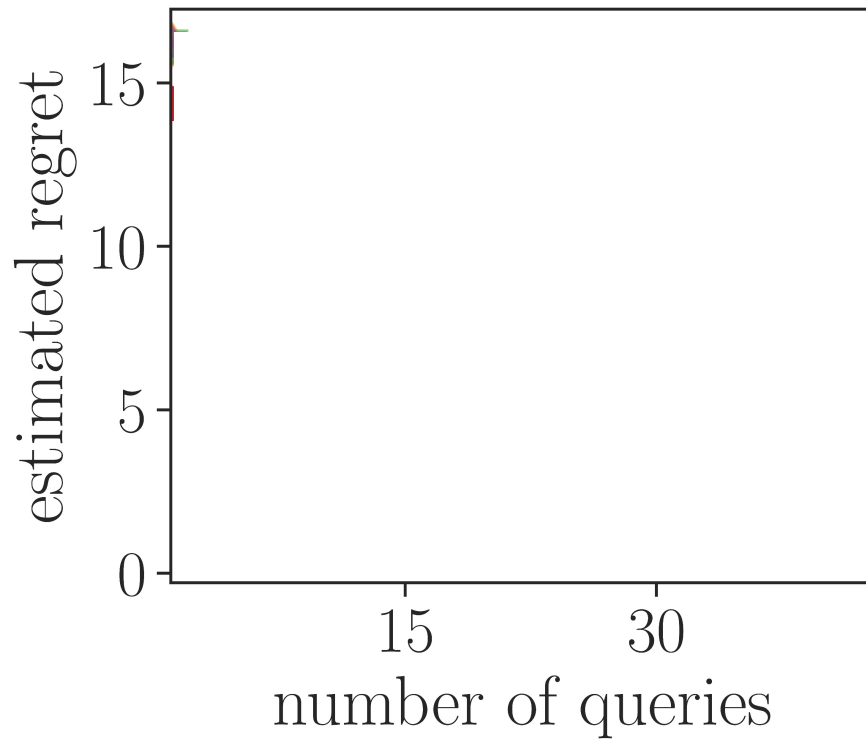# IDRL can learn a complex task in the MuJoCo simulator from numerical evaluations of clips of trajectories

GP reward model
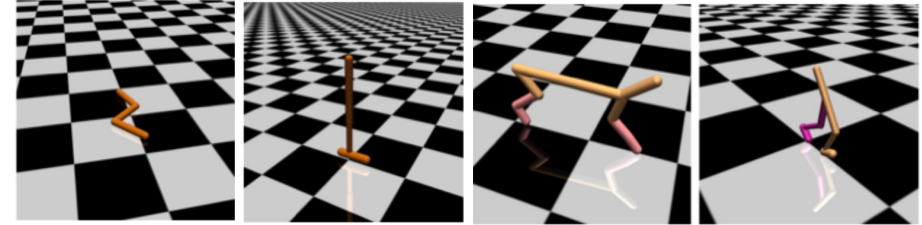


Information Directed Reward Learning (ours)

Information Gain on the Reward

Expected Improvement (EI)

Uniform Sampling

# IDRL can learn locomotion in MuJoCo from comparisons

- Normalized score averaged over multiple MuJoCo environments as a function of policy training steps
- Reward model trained from synthetic comparison queries
- Samples provided following a pre-defined schedule that is the same for all methods (more samples initially, less later)

Information Directed Reward Learning (ours)

IDRL w/o candidate policy rollouts

Information Gain on the Reward

Uniform Sampling

ETHzürich

# How to communicate task specifications to RL agents?



Task communication

Task representation

Reward function

Reward function

Explicit specification

Goal states

Constraints

Goal states

Constraints

How should we represent tasks that humans want RL agents to do?

Implicit specification

Demonstrations

Preference labels

RL / Policy optimization

Active learning

Which queries should RL agents make to learn what humans want most effectively?

# Many practical tasks naturally decompose into rewards and constraints

Robotics

"Pick up the screwdriver without hitting the wall."

"Cook me food but don't mess up the kitchen."

$r =$ „pick up screwdriver"
$c =$ „don't hit wall"

$r =$ „cook food"
$c =$ „keep kitchen clean"

Self-driving car

$r =$ „drive to grocery store"
$c =$ „obey driving rules"

"Get me to the grocery store as fast a possible."

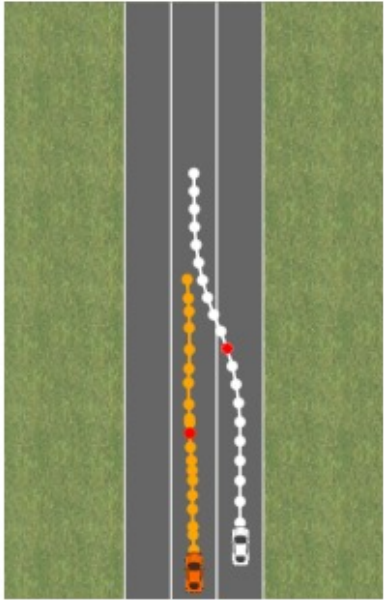# Constraint have can be more transferable and robust than rewards

Encode task as reward + penalty
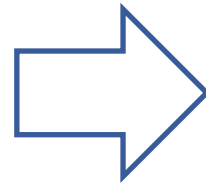


Base Scenario

Encode task as reward + constraint

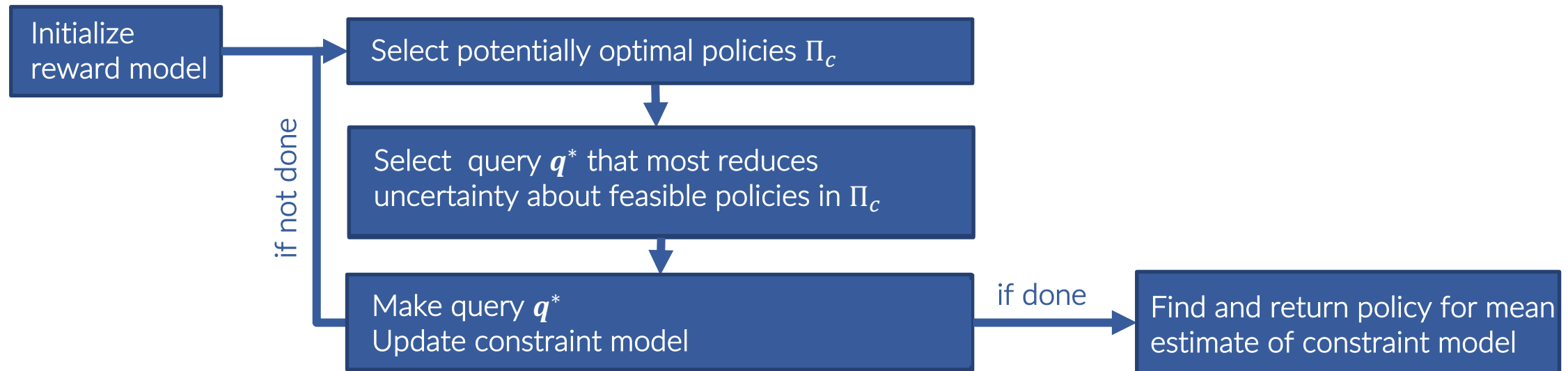# We can actively learn constraints similarly to learning rewards

**Information directed reward learning**

1. Which policies are potentially optimal?

2. How can we reduce uncertainty about their *optimality*?

**Adaptive constraint learning (ACOL)**
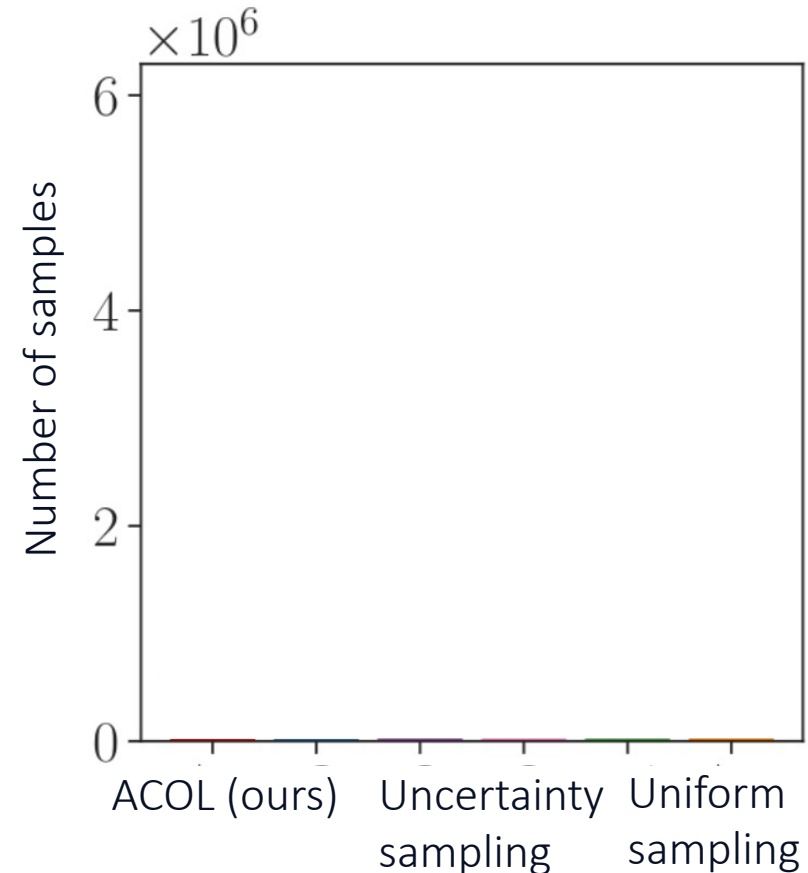
1. Which policies are potentially optimal?

2. How can we reduce uncertainty about their *feasibility*?

Initialize reward model

Select potentially optimal policies $\Pi_c$

Select query $q^*$ that most reduces uncertainty about feasible policies in $\Pi_c$

Make query $q^*$
Update constraint model

if not done

if done

Find and return policy for mean estimate of constraint model

# Active learning also improves sample efficiency for learning constraints

- We consider a driving environment with known reward but unknown constraints

- We can obtain binary samples wheather a trajectory is feasible or not

- How many samples to we need until we can identify the best constrained policy?

# Key takeaways

- For active reward learning, we should consider plausibly optimal policies
- Information directed reward learning (IDRL) provides a way to do so

- Rewards and constraints can be a **particularly robust** alternative to represent task specifications compared to reward functions only
- We can actively learn constraints with a similar method as IDRL

https://arxiv.org/abs/2102.12466

# Come talk to me if...

- ... you want to hear more technical details about anything I talked about.
- ... you work on an application that would benefit from learning task specifications from humans.
- ... you just want to chat.